

5 Some Ins and Outs with Python

Time (and the term) is a'wasting, so lets crank up your Canopy environment and get to work with some examples. In what follows it's assumed that you are using the Canopy "Editor," which includes a version of IPython, and so the command prompts will look like In '[n]:' without the single quotes. What "n" is depends on where you are in the current session. As I'm sure you realize by now since you have most likely use Canopy or IPython before, for many "Ins" in IPython there is an "Out" of some kind which is a result of the "In" that precedes it.

```
In [1]: outfile = open('myflatfile.txt','w') # open to write to a text file
```

In [2]: outfile # outfile is an open file 'object' in write mode. By default it's a text, not binary,
35 file:

```
In [3]: type(outfile)
```

Pardon the slight digression, above. Its purpose was to show what kind of Python 'object' outfile is. (Everything in Python is an object, right?) the .txt file name extension is optional. Now, let's create a text string and then write it to outfile:

```
In [4]: iLikeButter="""Slather me toast with a bargefull of butter,  
...:    and crown it with a bucket of Pythonberry jam."""
```

- 50 Note that those `"""` are three single quote characters. This is obviously a quote from a high cholesterol data science pirate. How many characters are in this string? Try the function `len(iLikeButter)`.

Write the string to outfile and then close outfile:

```
55 In [5]: outfile.write(iLikeButter)  
  
In [6]: outfile.close()      # it's good practice to close whatever you open
```

- 60 Now, where did you write this file to? It depends on your Canopy installation, what Canopy thinks your current working directory is, and some other things. One way to find out where it is is to look for it on your computer. Another way is to use the method for getting the current working directory in the `os` package:

```
65 In [7]: import os
```

```
In [8]: os.getcwd()
```

- 70 You'll get as a response a character string indicating what Canopy considers to be your working directory. A third way to find out what your current directory is from IPython is to type `!` (the “bang”), which escapes into your OS, followed by whatever command returns what the current directory is. On Linux and OS X, the command is `pwd`.

Let's read your file back in:

```
75 In [9]: infile=open('myflatfile.txt','r')      #read as a text file. For a binary, you'd use 'rb'  
  
In [10]: doYouWantButter=infile.read()  #reads the file contents into a string variable  
  
80 In [11]: doYouWantButter                #this should give you the iLikeButter string value  
  
In [12]: type(doYouWantButter)           # this should give you “str”  
  
In [13]: type(_)                         # should have the same result as above, right?  
85
```

When used as above, what does the underscore, `“_”`, represent?

- Next, let's read a text file with more than one line. The text file `louielouie.txt` has been provided to you. Pop it into the default directory for your session, the directory you identified
90 before. Then, open it for reading:

```
In [14]: kingsMenLouie=open('louielouie.txt','r')      #r' since this is a text file.
```

95 You could read this file with `kingsMenLouie.read()` to get a single string that you can print, examine, or manipulate, e.g.

```
In [15]: louielouie=kingsMenLouie.read()      # take a look at louielouie by typing its name
```

And, then you could split the lines in `louielouie` into a list of lines as strings:

100

```
In [15]: louielist=louielouie.split()        # lists are your Python Friend. (One of them, at least)
```

105 You could also read this file line by line with `readline()`. For example, to get the file contents into `louie` a string variable (and assuming that the file has been closed and opened again after the foregoing):

```
In [16]: louie="" # string var to hold the lines from the file. There are two " , """.
```

```
In [17] while True:
```

```
110 ...:     line=kingsMenLouie.readline()
...:     if not line:
...:         break
...:     louie+=line
...:
```

115

Give the above code a try to see what you get. Are `louie` and `louielouie` different? Try the command `louie==louielouie`.

120 Python usually does a good job closing files that have been opened, but it's good practice to do so explicitly whenever possible. This is especially true when you are writing data out to a file, as explicitly closing a file written to forces any remaining write operation to finish. Did you close all the files you opened, above?

125 A simple way to close a file you've written to is as follows. Suppose you want to write the character string `iLikeButter` to a file called `greaseitup.txt` in your current directory. If you do:

```
In [18] with open('greaseitup.txt','wt') as butterOut:
```

```
...:     butterOut.write(iLikeButter)
...:
```

130

the file will be closed automatically for you when your write operation is completed. Note the `'wt'` in the open statement. `'t'` is for text, but it's optional. if you include a `'b'` instead, you'll have a binary file instead of a text file.

135 The procedures for reading and writing binary files using `open`, `read`, etc. are for the most part the same as for text files, and so we're not going to spend time here on binary file input and output. We're shortly going to move on to reading and writing csv files, but before that let's take a look at the classic method for “serializing” (storing with permanency) python objects

called pickling.

A pickle file includes one or more Python objects that can be read back into Python that has a Python-specific, environment independent format. Two pickling packages in Python 2.7 are **pickle** and **cPickle**. **cPickle** is the faster of the two, and is the default algorithm in Python 3. So let's tinker with it here. First, import **cPickle**:

```
In [19]: import cPickle as pickle          # we'll just call it what it is, whether dill or sweet
```

Now let's pickle our *louielist* from above in a file in the current working directory.

```
In [20]: pickle.dump(louielist,open('louielist.p','wb'))
```

The above writes a binary pickle file. You can read the file back into Python like:

```
In [21]: louiesBack=pickle.load(open('louielist.p','rb'))
```

Did you get *louielist* back unchanged? Try `louielist == louiesBack` from the command prompt.

We're going to move ahead to consider csv files, but to do so we're going to make use of the very popular *pandas* package. So let's import *pandas* first, and then look at a simple example of a very useful *panda* object, the *DataFrame*.

```
In [22]: import pandas as pd              # "as" makes panda's nickname pd for this session
```

```
In [23]: import numpy as np              # numpy nicknamed np
```

```
In [24]: from pandas import DataFrame, Series    # for convenience
```

My guess is that you have used the **numpy** package before in your work or in a previous course. *DataFrame* and *Series* are very handy **pandas** data structures that can do yeoman work for you in your data mangement efforts.

By way of introduction, let's first read a little pickled **pandas** *DataFrame*. Put the file *littleDF2.p* in your default directory (or somewhere you can find it from in Canopy), and do:

```
In [25]: littleDF=pickle.load(open('littleDF2.p','rb'))    # if it's in you default dir
```

If you then ask *littleDF* to type itself out, like:

```
In [26]: littleDF
```

you'll get a pretty table-like graphical thing with labelled rows and columns. And if you do:

```
In [27]: print littleDF
```

185 You get:

Out [27]:

```
190          var1      var2      var3      var4
obs1 -1.228611  0.309656  0.963380  0.392899
obs2 -0.560940 -0.962692  0.042021 -0.289916
obs3  0.781060  1.140318  0.621084  0.682519
obs4 -0.092004 -1.178608  1.854705  1.011108
195 obs5  0.041939  1.290624 -0.313368 -0.749832
obs6 -0.156378  0.670368 -0.415693 -0.81301
```

littleDF is a **pandas** DataFrame:

200 In [28]: type(littleDF)

Out [28]: pandas.core.frame.DataFrame

205 A **pandas** DataFrame is a table-like data structure with columns that can be of different data types, and that has both row and column indices. The row index of littleDF is in the leftmost column, above, with index values of obs1 through obs6. The column index is across the top, var1 through var4. A DataFrame can have more complicated, hierarchical or “nested” indices.

210 A Series is like one column of a DataFrame. It's a kind of vector that has an associated index. A DataFrame can be thought of as a set of Series in the columns that share a single index, the row index.

littleDF is a **pandas** DataFrame, and DataFrames are Python objects. So, littleDF has attributes, and lots of them. To see them you can use the IPython tab completion feature:

215 In [29]: littleDF.<tab>

220 You'll see a long list of different attributes (methods and data) that littleDF has. A second <tab> should allow you to scroll through the attributes. Try littleDF.values to see the data, and littleDF.describe() to get summary statistics for the variables in the columns littleDF.columns.

DataFrame methods include different ways of converting to a different kind of data structure. For example, you can convert littleDF into a dict with :

225 In [30]: littleDF.to_dict('records')

or into csv:

230 In [31]: littleDF.to_csv()

or into several other formats. Note that all these methods have options set by arguments when calling them. To see what's available, do:

235 In [32]: littleDF.to_<tab>

Not surprisingly, you can convert other kinds of data structures into a DataFrame. Take a look:

240 In [33]: littleDF.from_<tab>)

Next, let's try writing out littleDF to a csv file called "littleDFcsv.csv:"

245 In [34]: csvOutFile=open('littleDFcsv.csv','wt') # this will be a text file

In [35]: csvOutFile.write(littleDF.to_csv())

In [36]: csvOutFile.close()

250 Now, take a look at the file you written using the editor of your choice. You can use the editor in Canopy, of course. Does your file have a header record? Are character values enclosed in quotes?

255 DataFrames and Series have many useful attributes and features, some of which we'll explore in upcoming exercises. But now let's try reading a less trivial csv file into a DataFrame. The file is xyzcust10.csv, and it should be available to you on Canvas. Take a look at it with your favorite text editor. Then, put it in a place you can find it from Canvas, and input it into a DataFrame:

260 In [37]: xyzcust10=pd.read_csv('xyzcust10.csv')

The file has 10 variables in it. The rows, or records, are XYZ customers. How many records are in xyzcust10?

265 What types of variables are in the columns of xyzcust10? To find out:

In [38]: xyzcust10.dtypes

Look at the first and last rows in xyzcust10:

270

In [39]: xyzcust10.head()

In [40]: xyzcust10.tail()

275 Note that in this file missing values for ZIP, ZIP4, and the nine digit ZIP are represented with zeros, "0's." The ZIPs could really also be coded as strings, rather than as integers, couldn't

they? Also, it looks like there might be two nine digit ZIP code variables. Are they the same? That is, are the values in these two variables the same for every row of data? How would you locate the rows in xyzcust10 that have a zero for ZIP or for ZIP4? We'll see in the next session's Python Practice.

Last but not least, be sure to pickle your xyzcust10 Data Frame so you can use it in the next Practice.